

Maximum entropy sampling in complex networks

Filippo Radicchi¹ and Claudio Castellano²

¹*Center for Complex Networks and Systems Research, School of Informatics and Computing,
Indiana University, Bloomington, Indiana 47408, USA**

²*Istituto dei Sistemi Complessi (ISC-CNR), Via dei Taurini 19, 00185 Roma, Italy*

Many real-world systems are characterized by stochastic dynamical rules where a complex network of dependencies among individual elements probabilistically determines their state. Even with full knowledge of the network structure and of the stochastic rules of the dynamical process, the ability to predict system configurations is generally characterized by large uncertainty. Sampling a fraction of the nodes and deterministically observing their state may help to reduce the uncertainty about the unobserved nodes. However, choosing these points of observation with the goal of maximizing predictive power is a highly nontrivial task, depending on the nature of the stochastic process and on the structure of the underlying network. Here, we introduce a computationally efficient algorithm to determine quasi-optimal solutions for arbitrary stochastic processes defined on generic sparse topologies. We show that the method is effective for various processes on different substrates. We further show how the method can be fruitfully used to identify the best nodes to label in semi-supervised probabilistic classification algorithms.

Stochastic phenomena are so ubiquitous in nature that they are studied in any field of science, including biology [1], ecology [2], physics [3], neuroscience [4], and finance [5]. Roughly speaking, in a stochastic system composed of multiple elements the states of the individual elements obey probabilistic rules that depend on the states of other elements in the system. In many cases of interest, the underlying structure that determines how elements depend one on the other is a sparse network [6]. To have a concrete example in mind, consider the seasonal flu spreading. The epidemics usually starts from a few epicenters. A person not immunized can contract the disease with a certain probability only if in contact with an infected individual. At the same time, infected people can recover and are thus no longer able to transmit the disease. The social network underlying the spreading process determines how the state of every individual depends on the state of the others. At any given time, the system is characterized by a certain degree of uncertainty, in the sense that different configurations have a non-vanishing probability to appear. Such an uncertainty is due to the stochastic nature of the spreading process, and it is present regardless of the level of knowledge that one has about the probabilistic rules of the process and about the contact pattern underlying the dynamics.

An obvious way to reduce uncertainty is to survey the system, i.e., observe the state of a sample of elements. In the example of flu spreading, this means obtaining full knowledge about the health state of some people. With such additional knowledge, the probabilistic prediction of the state of unobserved elements becomes less uncertain. In particular, the larger is the sample, the lower is the uncertainty, with the limiting case of null uncertainty when the entire system is observed. Resource constraints make however complete observation generally impossible: only a small fraction of elements of the system can be sampled. Is there an efficient way of identifying the best elements to observe in the sense that the uncertainty for the rest of system is minimized? The question is answered from an information theoretic point of view, by the so-called principle of maximum entropy sampling [7]. Its

rationale is very intuitive: to reduce the uncertainty about the system as much as possible, observations must be performed on the elements for which uncertainty is maximal.

Maximum entropy sampling is often contemplated in the literature as a possible solution to problems of experimental design [8]. An example is represented by the problem of choosing the placement of thermometers in a room to provide the most accurate picture of temperature in the entire room [9]. Some works, such as Refs. [10–12], have dealt with special settings where the solution to the problem of maximum entropy sampling can be efficiently approximated or achieved exactly with ad-hoc algorithms. These studies have, however, considered very small systems, where the scalability of the algorithms is not a concern. Further, the problem has been studied only in regular topologies, such as lattices or fully connected networks. No attention instead has been devoted so far to the maximum entropy sampling problem when the underlying topology is given by a large complex network structure. This is the focus of the present paper¹.

We consider an arbitrary stochastic process defined on a graph \mathbb{G} , composed of N nodes. Every node $i \in \mathbb{G}$ is characterized by a state variable x_i that can assume K distinct values. Connections among nodes in the network stand for dependency relations among their associated variables. Suppose the joint probability distribution $p(x_1, x_2, \dots, x_N)$ associated with each of the K^N possible configurations of the system is known. We are still left with potentially large uncertainty quantified by the information theoretic joint entropy of the graph

$$\mathcal{H}(\mathbb{G}) = - \sum_{x_1, x_2, \dots, x_N} p(x_1, x_2, \dots, x_N) \log_2[p(x_1, x_2, \dots, x_N)] . \quad (1)$$

The computation of $\mathcal{H}(\mathbb{G})$ involves a sum over K^N possible

¹ We stress that the problem discussed here is distinct from the problem of optimal sampling to reduce uncertainty on *topological* properties of a network [13].

configurations, thus becoming unfeasible even for small values of N .

Suppose we are allowed to observe a subset of nodes $\mathbb{O} \subseteq \mathbb{G}$. Observing these nodes means removing any uncertainty on their state, and thus considering the joint probability distribution of the unobserved part of the graph, $\mathbb{G} \setminus \mathbb{O}$, conditioned on the state of observed nodes \mathbb{O} , namely $p(x_{u_1}, \dots, x_{u_{N-O}} | x_{o_1}, \dots, x_{o_O}) = p(\mathbf{x}_{\mathbb{G} \setminus \mathbb{O}} | \mathbf{x}_{\mathbb{O}})$, where $u_1, \dots, u_{N-O} \in \mathbb{G} \setminus \mathbb{O}$, $o_1, \dots, o_O \in \mathbb{O}$, and we defined for shortness of notation $\mathbf{x}_{\mathbb{G} \setminus \mathbb{O}} = (x_{u_1}, \dots, x_{u_{N-O}})$ and $\mathbf{x}_{\mathbb{O}} = (x_{o_1}, \dots, x_{o_O})$. For a particular choice of the set \mathbb{O} , the uncertainty about the rest of the system is quantified by the conditional entropy

$$\mathcal{H}(\mathbb{G} \setminus \mathbb{O} | \mathbb{O}) = - \sum_{\mathbf{x}_{\mathbb{O}}} p(\mathbf{x}_{\mathbb{O}}) \sum_{\mathbf{x}_{\mathbb{G} \setminus \mathbb{O}}} p(\mathbf{x}_{\mathbb{G} \setminus \mathbb{O}} | \mathbf{x}_{\mathbb{O}}) \log_2 [p(\mathbf{x}_{\mathbb{G} \setminus \mathbb{O}} | \mathbf{x}_{\mathbb{O}})] . \quad (2)$$

For $\mathbb{O} = \emptyset$, Eq. (2) is identical to Eq. (1). For $\mathbb{O} = \mathbb{G}$, we have instead $\mathcal{H}(\mathbb{G} \setminus \mathbb{O} | \mathbb{O}) = \mathcal{H}(\emptyset) = 0$.

For a given number O of observed nodes, to reduce uncertainty on the system we want to observe the O nodes that minimize the conditional entropy of the unobserved part, Eq. (2). In particular, since $\mathcal{H}(\mathbb{G} \setminus \mathbb{O} | \mathbb{O}) = \mathcal{H}(\mathbb{G}) - \mathcal{H}(\mathbb{O})$, the minimization of Eq. (2) is equivalent to finding the group of nodes \mathbb{O}^* having maximum joint entropy, i.e.,

$$\mathbb{O}^* = \arg \max_{\mathbb{O}} \mathcal{H}(\mathbb{O}) , \quad (3)$$

where $\mathcal{H}(\mathbb{O}) = - \sum_{\mathbf{x}_{\mathbb{O}}} p(\mathbf{x}_{\mathbb{O}}) \log_2 [p(\mathbf{x}_{\mathbb{O}})]$. The maximization is performed over all sets \mathbb{O} of fixed size O . This principle is known as maximum entropy sampling, and the associated problem is NP -hard [7]. The exact solution of the optimization problem requires the consideration of all possible choices of the set \mathbb{O} , and for each of them the computation of the associated joint entropy. The computational complexity of both operations scales exponentially with O .

A quasi-optimal solution to the problem can however be obtained at a reduced computational cost, exploiting the submodularity of the entropy function [14]. Such a property allows to implement a greedy strategy where the set of observed nodes is constructed in a sequential manner, leading to a solution that is provably close to the optimum. It provides a solution corresponding to a value of the function to be optimized that is at least $(1 - 1/e) = 0.66 \dots$ times the value of the global maximum. In the context of our problem, the greedy algorithm consists in sequentially adding, to the set of observed nodes, the node with maximal entropy conditioned to the set of variables that are already observed. More specifically, the algorithm starts at stage $t = 0$ with an empty set of observed nodes, $\mathbb{O}_{t=0} = \emptyset$. The t -th point of observation, namely o_t , is chosen, among the nodes not yet part of the set of observed nodes $\mathbb{O}_{t-1} = \{o_1, o_2, \dots, o_{t-1}\}$, according to the rule

$$o_t = \arg \max_{i \notin \mathbb{O}_{t-1}} \mathcal{H}(i | o_1, \dots, o_{t-1}) . \quad (4)$$

The algorithm can be run up to arbitrary values $1 \leq t \leq N$. This procedure addresses the issue of the extensive search

over all possible groups of nodes. However, at every stage t , the computation of each of the $N - (t - 1)$ conditional entropies in Eq. (4) still requires a number of operations scaling as K^t . This makes the algorithm usable only for the construction of very small sets of observed nodes.

The issue of scalability cannot be solved in general. However, if the underlying graph \mathbb{G} is a tree, a computationally efficient algorithm to compute the conditional entropies in Eq. (4) can be deployed as follows. This is the main contribution of the present paper. We suppose that the tree is connected, so that every node is reachable from every other node along one and only one path. If the graph is composed of multiple connected components, these can be considered independently. The algorithm makes use of three main properties of conditional entropy: (i) chain rule, (ii) conditional independence, and (iii) Bayes rule. Suppose we are at stage t of the greedy algorithm. We need to compute the conditional entropy $\mathcal{H}(i | o_1, \dots, o_{t-1})$ for a generic node $i \notin \mathbb{O}_{t-1}$. Thanks to the Bayes rule, we can write

$$\mathcal{H}(i | o_1, \dots, o_{t-1}) = \mathcal{H}(i) + \mathcal{H}(o_1, \dots, o_{t-1} | i) - \mathcal{H}(o_1, \dots, o_{t-1}) . \quad (5)$$

In the above expression, $\mathcal{H}(i)$ is the unconditional entropy of node i , $\mathcal{H}(o_1, \dots, o_{t-1} | i)$ is the joint entropy of the nodes o_1, \dots, o_{t-1} conditioned to node i , and $\mathcal{H}(o_1, \dots, o_{t-1})$ is the joint unconditional entropy of the nodes o_1, \dots, o_{t-1} . The first term on the r.h.s. of Eq. (5) can be easily and cheaply estimated.

Given that the network is a tree, the second term can be written (see Appendix for details), as the sum of pairwise conditional entropies, one per observed node

$$\mathcal{H}(o_1, \dots, o_{t-1} | i) = \sum_{j=1}^{t-1} \mathcal{H}(o_j | s_{o_j}) . \quad (6)$$

Every observed node corresponds to a term in the sum, given by the entropy associated with that observed node, o_j , conditioned to another node $s_{o_j} \in \{i\} \cup (\mathbb{O}_{t-1} \setminus \{o_j\})$. Such a node s_{o_j} is either the first observed node encountered along the unique path connecting o_j to i , or node i itself.

Finally, thanks to the chain rule, for the rightmost term in Eq. (5) we know that

$$\mathcal{H}(o_1, \dots, o_{t-1}) = \mathcal{H}(o_{t-1} | o_1, \dots, o_{t-2}) + \mathcal{H}(o_1, \dots, o_{t-2}) . \quad (7)$$

Hence, the rightmost term of Eq. (5) can be computed at each stage t by reusing results of the algorithm in previous stages. We note that the entire procedure can be used for arbitrary sequences of added nodes, not just the one decided according to the greedy algorithm. Hence, Eqs. (5-7) allow us to compute the joint entropy of any subset of nodes, including the entire graph. This is a side product of the algorithm, yet a very important and useful result.

The algorithm requires prior knowledge of the unconditional entropy of the individual nodes, and the pairwise entropy among pairs of nodes. These quantities can be esti-

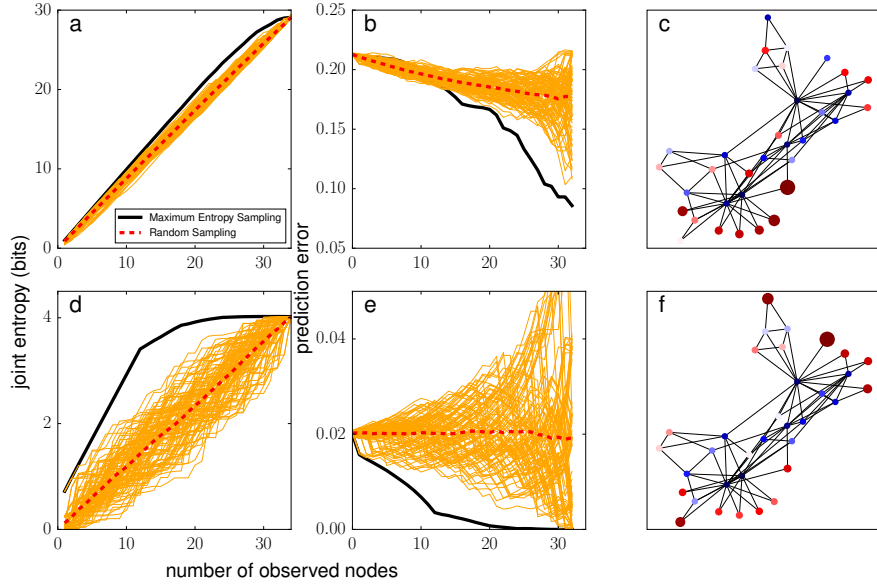


Figure 1. Maximum entropy sampling in real networks. We study the bond percolation model applied to the Zachary karate club network for two different values of the occupation probability p . In panels a, b and c, $p = 0.3$. a) The plot shows how the joint entropy of the set of observed nodes behaves as a function of its size. Two sampling techniques are considered here: the method introduced in this paper, i.e., maximum entropy sampling (black full line); ii) random sampling, where nodes are introduced in the set of observed nodes in random order (thin orange lines are 100 single realizations, red dashed line stands for their average value). b) Quadratic error $\epsilon^2 = \frac{1}{N-O} \sum_{i \notin O} [p(x_i = 1 | \mathbf{x}_O = \tilde{\mathbf{x}}_O) - \tilde{x}_i]^2$ in the prediction of the states of the unobserved nodes as a function of the size O of the set of observed nodes. The inference test is based on $T = 10000$ independent simulations of the bond percolation process. c) Graph visualization of the network. Color (dark red to dark blue) and size (large to small) of the nodes stand for the order in which they are introduced in the set of observed nodes according to maximum entropy sampling. Panels d, e and f are identical to panels a, b, and c, respectively, but are obtained for $p = 0.8$.

mated either from numerical simulations or experimental observations of the stochastic process, or theoretically, using for example a belief-propagation algorithm [15]. It is worth to remark the great advantages brought by our approach compared with the naive method to compute joint conditional entropies when a system is studied with numerical simulations (or experimental observations). This may be the case of most situations. If the size of the set of observed nodes is O , the maximum value of the entropy is $O \log_2 K$. This corresponds to the case in which each of the K^O configurations has exactly the same probability to occur. If the number of simulations is T , the maximum value of the entropy that can be measured with the naive approach is $\log_2 T$. For large values of O , the number of simulations will be unavoidably $T \ll K^O$, thus leading to systematically biased estimates. With our method instead, the maximum value of the pairwise conditional entropy between a pair of nodes is $\log_2 T$. The total entropy of the set of observed nodes will be given by O terms of this type, leading to a maximum possible value equal to $O \log_2 T$. A number of simulations $T \gg K$, computationally feasible in most situations, is then sufficient to avoid numerical problems.

From a computational point of view, the running time scales as N^3 in the worst-case scenario. We found, however, that some computational tricks allow for a great reduction of the complexity [14, 16], which effectively scales as $N^2 \log(N)$ (see Appendix). This makes the present algorithm easily ap-

plicable even to large systems.

The above algorithm is exact if the graph G is a tree. We argue that the algorithm can be still effectively used in loopy but sparse graphs, as many of the networks describing real systems [6]. In the case of loopy graphs, the conditional entropy $\mathcal{H}(o_1, \dots, o_{t-1} | i)$ is still computed under the tree assumption by generating a spanning tree rooted in i . This provides us with an upper-bound of its true value. The rooted tree can be generated arbitrarily. However, to keep the upper-bound as tight as possible, we use a Dijkstra-like algorithm suitably modified for this context (see Appendix). The results presented below are based on this choice. We believe that our algorithm may be useful for any stochastic process where the dependence among variables is represented by a sparse network. This is a common setting in several realistic scenarios. We expect a decrease of performance not only in dense networks, but also in sparse networks with spatial embedding (as for example road networks or power grids), or networks not compatible with the locally tree-like ansatz (as for example collaboration networks). In practical applications, assuming that the network structure is perfectly known may not be necessarily true. The method can be modified to include such an additional degree of uncertainty, as long as sparsity of the network is among the ingredients of the stochastic model.

As a proof of concept to verify the effectiveness of the algorithm, we consider the bond percolation model [17]. This

is a prototypical example of a stochastic process on a complex network topology with relevance in the analysis of networks robustness [18], and of spreading phenomena such as those belonging to the Susceptible-Infected-Recovered class [19]. In a single realization of the model, every edge is considered active or occupied with bond occupation probability p . Nodes that are connected by at least a path formed by occupied edges form clusters. We focus our attention on the largest among these clusters. We suppose that the state of a given node i can take only $K = 2$ values: $x_i = 1$ if node i is part of the largest cluster in the network, and $x_i = 0$, otherwise. For simplicity, we assume that the value of the bond occupation probability p is known in advance. To estimate the unconditional entropy $\mathcal{H}(i)$ of a generic node i , and the pairwise conditional entropy $\mathcal{H}(j|i)$ of a generic pair of nodes i and j , we rely on the belief-propagation algorithm recently introduced in Ref. [20]. This algorithm provides us with a set of probabilities $p(x_i = 1)$ for every node i to be part of the largest cluster. To obtain conditional probabilities we can rely on the same algorithm by simply blocking the value of the conditional variables. The belief-propagation algorithm is developed under the treelike approximation, so it may not provide the best estimate of the marginal and conditional entropy in loopy networks [21]. Further, we note that the belief-propagation algorithm provides a precise value of the percolation threshold, namely p_c , below which the largest cluster doesn't exist although this is not true for a network with finite size. The consequence of this fact is that, for $0 \leq p \leq p_c$, the entropy of the system is estimated to be equal to zero. Null entropy is also associated to the system for $p = 1$. At intermediate values of p , the system is characterized by a non-null uncertainty. This is the regime where a proper sampling of the network is needed. In our analysis, we compare the performance of maximum entropy sampling with random sampling, i.e., a simple strategy where the set of observed nodes is iteratively constructed by adding nodes in random order. To compare the performance of two sampling strategies, we use two independent tests. The first consists in the comparison between the values of the joint entropy of the set observed nodes. The higher this quantity, the better the strategy. The second test is based on inference. We simulate the bond percolation model T independent times. In every simulation, we identify the largest cluster (if more clusters have maximal size, we randomly choose one of them to be the largest), and we assign a value $\tilde{x}_i = 1$ to every node i belonging to the largest cluster, or $\tilde{x}_i = 0$, otherwise. Then, for a given set of observed nodes \mathcal{O} , we use the belief-propagation algorithm by blocking the value of the variables obtained in the simulation for all observed nodes, that means $p(x_i = 1) = \tilde{x}_i$ for all $i \in \mathcal{O}$. The algorithm provides us with the marginal conditional probability $p(x_i = 1 | x_{o_1} = \tilde{x}_{o_1}, \dots, x_{o_o} = \tilde{x}_{o_o}) = p(x_i = 1 | \mathbf{x}_\mathcal{O} = \tilde{\mathbf{x}}_\mathcal{O})$ to be part of the largest cluster for every node $i \notin \mathcal{O}$. We note that this doesn't exactly correspond to make inference on the unobserved part of the networks, as the required probability would have been $p(\mathbf{x}_{G \setminus \mathcal{O}} | \mathbf{x}_\mathcal{O} = \tilde{\mathbf{x}}_\mathcal{O})$. We finally compare the probabilistic prediction with the actual configuration. Fig-

ure 1 summarizes the results obtained for the bond percolation model applied to the Zachary karate club network [22]. The percolation threshold of the network is $p_c = 0.1889$ according to the belief-propagation method [23]. For values of p close to p_c , the uncertainty of the system is high, as it appears from the joint entropy in Figure 1a. Sampling the network according to the maximum entropy principle allows us to systematically construct a set of observed nodes with joint entropy that is slightly larger than the one obtained for a set of observed nodes chosen at random. The prediction error obtained with both strategies is systematically decreasing as the number of observed nodes increases (Fig. 1b), with maximum entropy sampling obtaining slightly better performances. We note that maximum entropy sampling tends to first select peripheral and low-degree nodes, while central and large-degree nodes are generally included in the set of observed nodes only at the end of the sampling procedure (Fig. 1c). For large value of the occupation probability p , the difference in performance between maximum entropy and random sampling becomes more evident. This is visible in both Figures 1d and 1e. As soon as we reach a set of observed nodes that has joint entropy equal to the one of the entire graph (around $O \simeq 20$), the rest of the system becomes essentially deterministic, and the state of the unobserved part of the network is predicted with no error. Still, best observation points are given by low-degree nodes, however different from those identified at low p values (Fig. 1f). Random sampling has instead very poor performances. This is particularly apparent from the large values of the prediction error. In a single realization of random sampling, the fact that the prediction error may increase as the size of observed set increases is due to the fact that the belief-propagation algorithm may provide "wrong" predictions in a loopy network. The main message of the analysis is, however, particularly important: with random sampling there is a concrete risk of wasting resources, as we do not learn much about the rest of the system, even by observing a large fraction of its nodes.

From the former analysis, it seems that peripheral low-degree nodes are the best points of observation to infer the state of the rest of the system. This is somehow reminiscent of what noticed in Refs. [9, 12], although the systems studied there were defined on two-dimensional lattices. While finding the best points of observation at the border of the system could be a rather general trend in several stochastic systems (and may be fruitfully used in heuristic methods for the optimal placement of sensors in a system), we argue that this is not always true. Several systems can exhibit different behaviors and, even for the same system and the same topology, the behaviour can be modified by particular choices of the system parameters. In Figure A3, we show results obtained from the analysis of the bond percolation model on an uncorrelated configuration model with scale-free degree distribution [26]. We see that top nodes in the set of observed nodes are chosen differently depending on the value of the bond occupation probability p . To strengthen this message, in Figure A4, we consider the Susceptible-Infected-Susceptible model relying on the approach of Ref. [27]. Qualitative results are identical

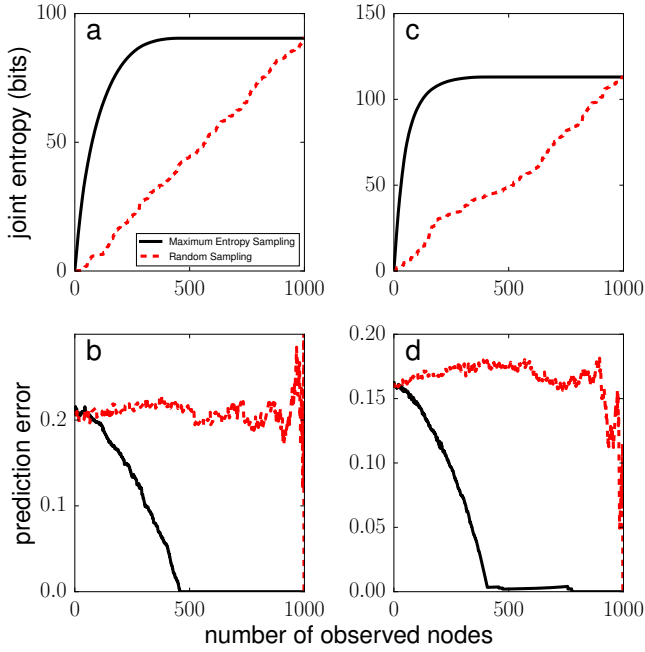


Figure 2. Maximum entropy sampling in semi-supervised community detection. We consider graphs with $N = 1000$ nodes divided in $K = 2$ (panels a and b) and $K = 4$ (panels c and d) groups. Average internal degree is $\langle k_{in} \rangle = 10$. 90% of the nodes have average external degree $\langle k_{ext} \rangle = 10$. For the remaining 10%, we set $k_{ext} = 0$. We use the expectation maximization algorithm of Ref. [24] to recover communities. The starting configuration of the iterative algorithm is the true community assignment of the nodes as suggested in Ref. [25]. a and c) Joint entropy of the set of observed nodes as a function of the number of observed nodes. Maximum entropy sampling (black full line) is compared against a single instance of random sampling (red dashed line). b and d) Error in the prediction of the labels of the unobserved part of the system as a function of the number of observed nodes. The error is calculated according to the formula $\epsilon^2 = \frac{1}{(N-O)K} \sum_{i \notin O} \sum_{x_i=1}^K [\delta_{x_i, \tilde{x}_i} - p(x_i | \mathbf{x}_O = \tilde{\mathbf{x}}_O)]^2$, where \tilde{x}_i is the index of the pre-assigned community of node i , and $\delta_{x,y} = 1$ if $x = y$, whereas $\delta_{x,y} = 0$, otherwise.

to those of the percolation model. Further in the Appendix, we study analytically the behavior of the Independent Cascade Model and show that, in star-like networks, the choice of the best nodes to observe is highly sensible to the parameter of the model, and the initial configuration used for the stochastic dynamical process.

So far, we focused on stochastic processes taking place on networks, but our method may be important also in classification problems on networks, such as those considered by probabilistic semi-supervised machine learning methods [15]. In this context, the set of observed nodes represents the set of points to be labelled. Among the applications of semi-supervised classification algorithms in sparse networks is community detection, where many probabilistic models for the detection of clusters have been considered [28]. We conclude the paper with an analysis in this direction (Fig. 2). We

consider artificial networks introduced in Ref. [29], where N nodes are divided in K groups of identical size. Every node has an internal degree coming from a Poisson distribution with average $\langle k_{in} \rangle$, and an external degree towards any other group coming from a Poisson distribution with average degree $\langle k_{ext} \rangle$ (average external degree to all other groups is thus $(K - 1)\langle k_{ext} \rangle$). To create imbalance among nodes in the groups, we force external degrees of a fraction $\alpha = 0.1$ of nodes to be equal to zero, so that these nodes have only connections to nodes within their pre-assigned groups. For given internal and external degree sequences, the graph is created according to the same scheme of the configuration model. We apply the expectation maximization algorithm of Ref. [24] to identify communities. The algorithm provides us with a probability $p(x_i)$, with $x_i = 1, \dots, K$, for every node i to be in a group. The algorithm provides us also with conditional probabilities obtained by blocking the value of some nodes. We use these quantities to compute marginal and pairwise conditional entropies required by our method, and to perform inference by blocking the values of the nodes in the observed set. Results of the analysis are reported in Figure 2. Maximum entropy sampling allows us to reduce uncertainty in a systematic manner. If the community detection method has sufficient knowledge, then it is even able to recover perfectly the labels of the unobserved nodes. Random sampling instead doesn't help much, as the community detection algorithm learns very slowly if labeled nodes are chosen without a precise criterion. Community detection is just a particular case study. So far, the maximum entropy principle has been considered only in some semi-supervised algorithms that incorporate active learning [30]. The approach has been never considered in the context of sparse network structures. We believe that the results of this paper may be helpful for future research in machine learning problems on sparse networks.

We thank A. Fagheh, A. Flammini, and S. Fortunato for comments on the manuscript. F.R. acknowledges support from the National Science Foundation (Grant No. CMMI-1552487) and the U.S. Army Research Office (Grant No. W911NF-16-1-0104).

* filiradi@indiana.edu

- [1] Paul C Bressloff, *Stochastic processes in cell biology*, Vol. 41 (Springer, 2014).
- [2] Russell Lande, Steinar Engen, and Bernt-Erik Saether, *Stochastic population dynamics in ecology and conservation* (Oxford University Press on Demand, 2003).
- [3] N Gr van Kampen, *Stochastic processes in physics and chemistry* (Elsevier, 1995).
- [4] Carlo Laing and Gabriel J Lord, *Stochastic methods in neuroscience* (Oxford University Press, 2010).
- [5] Frederi Viens, *Stochastic Processes: From Physics to Finance* (Taylor & Francis, 2002).
- [6] Mark Newman, *Networks: an introduction* (Oxford University Press, 2010).
- [7] Michael C Shewry and Henry P Wynn, "Maximum entropy

- sampling,” *Journal of Applied Statistics* **14**, 165–170 (1987).
- [8] Kathryn Chaloner and Isabella Verdinelli, “Bayesian experimental design: A review,” *Statistical Science*, 273–304 (1995).
 - [9] Andreas Krause, Ajit Singh, and Carlos Guestrin, “Near-optimal sensor placements in gaussian processes: Theory, efficient algorithms and empirical studies,” *Journal of Machine Learning Research* **9**, 235–284 (2008).
 - [10] Chun-Wa Ko, Jon Lee, and Maurice Queyranne, “An exact algorithm for maximum entropy sampling,” *Operations Research* **43**, 684–691 (1995).
 - [11] Jon Lee, “Maximum entropy sampling,” *Encyclopedia of Environmetrics* (2002).
 - [12] Carlos Guestrin, Andreas Krause, and Ajit Paul Singh, “Near-optimal sensor placements in gaussian processes,” in *Proceedings of the 22nd international conference on Machine learning* (ACM, 2005) pp. 265–272.
 - [13] Jure Leskovec and Christos Faloutsos, “Sampling from large graphs,” in *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’06 (ACM, New York, NY, USA, 2006) pp. 631–636.
 - [14] Andreas Krause and Daniel Golovin, “Submodular function maximization,” *Tractability: Practical Approaches to Hard Problems* **3**, 8 (2012).
 - [15] Christopher M Bishop, “Pattern recognition,” *Machine Learning* **128**, 1–58 (2006).
 - [16] Jure Leskovec, Andreas Krause, Carlos Guestrin, Christos Faloutsos, Jeanne VanBriesen, and Natalie Glance, “Cost-effective outbreak detection in networks,” in *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining* (ACM, 2007) pp. 420–429.
 - [17] Dietrich Stauffer and Ammon Aharony, *Introduction to percolation theory* (CRC press, 1994).
 - [18] Réka Albert, Hawoong Jeong, and Albert-László Barabási, “Error and attack tolerance of complex networks,” *Nature* **406**, 378–382 (2000).
 - [19] Romualdo Pastor-Satorras, Claudio Castellano, Piet Van Mieghem, and Alessandro Vespignani, “Epidemic processes in complex networks,” *Reviews of Modern Physics* **87**, 925–979 (2015).
 - [20] Brian Karrer, Mark EJ Newman, and Lenka Zdeborová, “Percolation on sparse networks,” *Physical Review Letters* **113**, 208702 (2014).
 - [21] Filippo Radicchi and Claudio Castellano, “Beyond the locally treelike approximation for percolation on real networks,” *Physical Review E* **93**, 030302 (2016).
 - [22] Wayne W. Zachary, “An information flow model for conflict and fission in small groups,” *Journal of Anthropological Research* **33**, 452–473 (1977).
 - [23] Filippo Radicchi, “Predicting percolation thresholds in networks,” *Physical Review E* **91**, 010801 (2015).
 - [24] Mark EJ Newman and Elizabeth A Leicht, “Mixture models and exploratory analysis in networks,” *Proceedings of the National Academy of Sciences USA* **104**, 9564–9569 (2007).
 - [25] Andrea Lancichinetti and Santo Fortunato, “Community detection algorithms: a comparative analysis,” *Physical review E* **80**, 056117 (2009).
 - [26] Michele Catanzaro, Marián Boguñá, and Romualdo Pastor-Satorras, “Generation of uncorrelated random scale-free networks,” *Physical Review E* **71**, 027103 (2005).
 - [27] A. V. Goltsev, S. N. Dorogovtsev, J. G. Oliveira, and J. F. F. Mendes, “Localization and spreading of diseases in complex networks,” *Physical Review Letters* **109**, 128702 (2012).
 - [28] Santo Fortunato, “Community detection in graphs,” *Physics Reports* **486**, 75–174 (2010).
 - [29] Filippo Radicchi, “A paradox in community detection,” *EPL (Europhysics Letters)* **106**, 38001 (2014).
 - [30] Xiaojin Zhu, John Lafferty, and Ronald Rosenfeld, *Semi-supervised learning with graphs* (Carnegie Mellon University, language technologies institute, school of computer science, 2005).

APPENDIX

Reduction of joint conditional entropy in Eq. (5) to a sum of pairwise conditional entropies

The second term on the r.h.s. of Eq. (5) can be greatly simplified when the interaction pattern is a tree. We show this in detail by considering the structure depicted in Fig. A1. In that case the second term of Eq. (5) is $\mathcal{H}(o_j, o_w, o_s, o_r|i)$. By using the chain

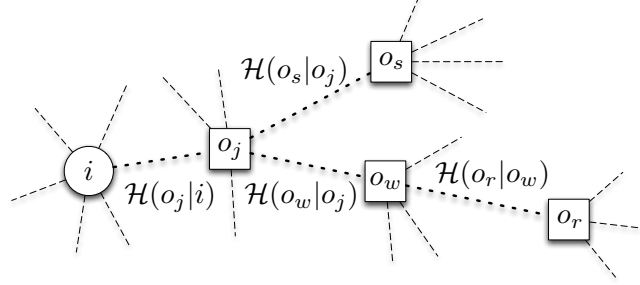


Figure A1. In a tree, the joint entropy of the observed nodes o_j , o_s , o_w , and o_r conditioned to node i can be broken in the sum of pairwise conditional entropies. Each observed node is responsible for one of these contributions, in the sense that the variable associated with it is conditioned to a variable associated with another node. The node that acts on the conditional part of each pairwise entropy is given either by the first observed node encountered along the path towards i or by node i itself. In the illustration above, the path connecting o_r to i passes through node o_w , giving rise to the contribution $\mathcal{H}(o_r|o_w)$. The paths connecting o_s and o_w to i pass both through node o_j , giving rise to the contributions $\mathcal{H}(o_s|o_j)$ and $\mathcal{H}(o_w|o_j)$, respectively. There are instead no observed nodes along the path between node o_j and node i , thus the contribution of node o_j to the joint entropy is $\mathcal{H}(o_j|i)$.

rule it can be written as

$$\mathcal{H}(o_j, o_w, o_s, o_r|i) = \mathcal{H}(o_j|i) + \mathcal{H}(o_w, o_s, o_r|o_j, i) = \mathcal{H}(o_j|i) + \mathcal{H}(o_w, o_s, o_r|o_j) \quad (\text{A1})$$

where the last step takes into account the tree topology. Again the consideration of the topology indicates that of o_s , conditioned to the value of o_j is independent from the values of o_w and o_r conditioned to o_j (conditional independence). Hence

$$\mathcal{H}(o_j, o_w, o_s, o_r|i) = \mathcal{H}(o_j|i) + \mathcal{H}(o_s|o_j) + \mathcal{H}(o_w, o_r|o_j) \quad (\text{A2})$$

The application of the chain rule on the last term leads to

$$\mathcal{H}(o_j, o_w, o_s, o_r|i) = \mathcal{H}(o_j|i) + \mathcal{H}(o_s|o_j) + \mathcal{H}(o_w|o_j) + \mathcal{H}(o_r|o_w) \quad (\text{A3})$$

The total conditional entropy is thus reduced to the sum of pairwise conditional entropies, one for each observed node.

The argument can be readily generalized to any tree, by multiple applications of the chain rule and of conditional independence, leading to

$$\mathcal{H}(o_1, \dots, o_{t-1}|i) = \sum_{j=1}^{t-1} \mathcal{H}(o_j|s_{o_j}). \quad (\text{A4})$$

In the sum, the variable associated with each observed node o_j is conditioned to a variable associated with another node $s_{o_j} \in \{i\} \cup (\mathcal{O}_{t-1} \setminus \{o_j\})$. The node s_{o_j} that acts on the conditional part of each pairwise entropy is given either by the first observed node encountered along the path towards i or by node i itself. Notice that dashed lines in Fig. A1 need not to be direct connections: other unobserved nodes may lie between observed ones.

Computational Complexity of the Algorithm

Our algorithm requires us to have prior knowledge of the unconditional entropy of every single node in the graph. The computation of these quantities scales as KN . The algorithm also requires prior knowledge of the pairwise conditional entropy among all pairs of nodes, whose computation scales as K^2N^2 . At stage t of the algorithm, we need to find the node i that

maximizes the entropy $\mathcal{H}(i|o_1, \dots, o_{t-1})$. This means that we have to compute the function for every unobserved node. We thus have to walk back on the tree from any already observed node towards the node whose entropy is being computed, visiting all edges M of the tree, with an algorithm that scales as M . As we need to perform such an operation for all nodes that are still in the set of unobserved nodes, the computational complexity of estimating the conditional entropy of all unobserved nodes is NM . These operations must be repeated at any stage of the greedy algorithm adding a factor N to the computational complexity. In conclusion, the computational complexity of the entire algorithm is $N^2 M \sim N^3$. Several considerations and computational techniques may be used to speed up the algorithm. First, knowledge of the pairwise entropy is not required for all pairs of nodes. *A priori*, we don't know which pairwise entropies will be used by the algorithm, but we can compute them on-the-fly when they are needed. Second, we know that the entropy of individual nodes conditioned to the set of observed nodes can only decrease during the algorithm, i.e., $\mathcal{H}(i|o_1, \dots, o_{t-2}) \geq \mathcal{H}(i|o_1, \dots, o_{t-2}, o_{t-1})$. We can therefore use a lazy algorithm that computes Eq. (5) for an unobserved node only if needed [14, 16]. To keep track of the ranking of unobserved nodes in a computationally cheap way, we can make use of a standard queue algorithm. Using these tricks, the effective computational complexity of the algorithm is greatly reduced (Fig. A2).

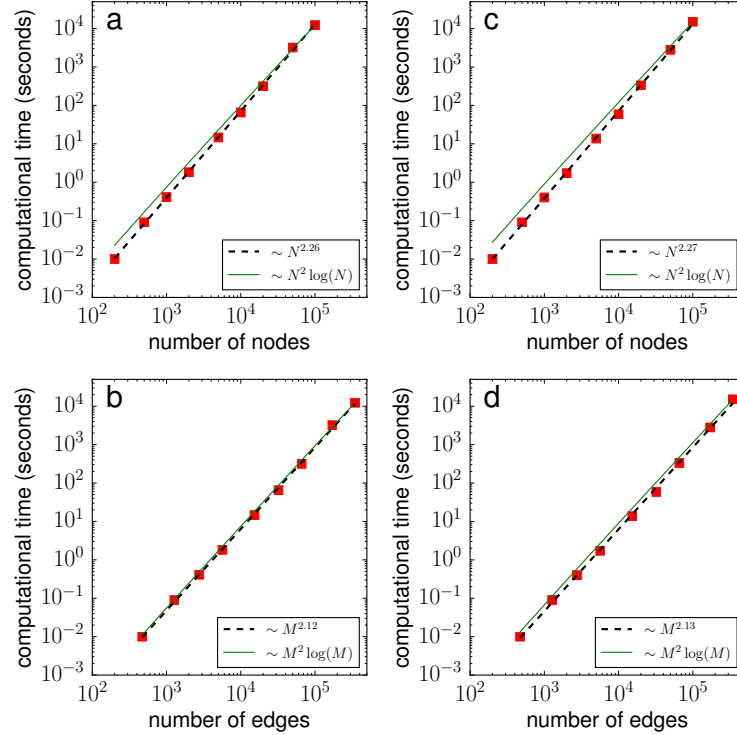


Figure A2. Computational complexity of the maximum entropy sampling algorithm. We measure the total time T required by the algorithm to generate sets of observed nodes from size 1 to N . Computations relative to marginal and conditional entropies are not included in our measure. As a representative case we consider here bond percolation applied to scale-free networks. The degree distribution of the networks is $P(k) \sim k^{-5/2}$ for $k \in [3, \sqrt{N}]$, and $P(k) = 0$. We further consider different values of the occupation probability p . a) Computational time as a function of the network size for $p = 0.5$. The dashed black line is a power-law fit $T \sim N^\beta$, with $\beta = 2.26$. The full green line corresponds instead to a fit $T \sim N^2 \log(N)$. b) Same as in panel a but showing the dependence of the computational time from the total number edges in the network. c) Same as in panel a, but for $p = 2p_c$. The critical values of the occupation probability of each network are computed according to the method of Ref. [20]. d) Same as in panel c, but showing the dependence of the computational time from the total number edges in the network.

A Dijkstra-like algorithm to determine the tree with minimal entropy rooted in a specific node

The following provides a description of the algorithm used to determine the tree with minimal entropy rooted in a specific node. The algorithm is a variant of the well-known Dijkstra's algorithm. We used this algorithm whenever we wanted to provide an estimate of the conditional entropy $\mathcal{H}(o_1, \dots, o_{t-1}|i)$ in loopy networks. Here, i is the index of node that acts as a root of the tree, o_1, o_2, \dots, o_{t-1} are instead the indices of the nodes already observed. We assume that these variables are given, as well as pairwise conditional entropies among pairs of nodes, and the entire topology of the network.

Define \vec{v} , with $v_n = -1$ for all $n = 1, \dots, N$. This vector serves to account for the fact that a node has been already visited. Define \vec{s} , with $s_n = -1$ for all $n = 1, \dots, N$, except for $s_i = i$. This vector serves to keep track of the pairwise conditional entropies that enter in the sum. Define \vec{o} , where $o_n = -1$ for all $n = 1, \dots, N$, except for the observed nodes $o_{o_1} = o_{o_2} = \dots = o_{o_{t-1}} = 1$. Finally, define the vector \vec{d} , with $d_n = \infty$ for all $n = 1, \dots, N$, except for $d_i = 0$. This vector represents the distance of a generic node from node i . Distance is measured in terms of the value of pairwise entropies along the path.

1. Select $n = \arg \min_{m|v_m < 0} d_m$, i.e., the node at minimal distance from i among those not yet visited.
2. Look at all neighbors of node n . For a given neighbor m that has not yet been visited, i.e., $v_m < 0$, apply one of the following mutually exclusive operations:
 - a) if $o_m > 0$ and $s_m < 0$, set $d_m = d_n + \mathcal{H}(m|s_n)$. If $o_n > 0$ or $n == i$, set $s_m = n$. Otherwise, set $s_m = s_n$. This part applies to observed nodes that we didn't have yet considered.
 - b) if $o_m > 0$ and $s_m > 0$, compute $\delta = \mathcal{H}(m|s_n) - \mathcal{H}(m|s_m)$. If $\delta < 0$, then set $s_m = s_n$ and $d_m = d_n + \mathcal{H}(m|s_n)$. This part applies to observed nodes that we have already considered, but for which we may find a shorter path towards i .
 - c) if $o_m < 0$, set $s_m = s_n$ and $d_m = d_n$. This part applies to non observed nodes.
3. Set $v_n = 1$, and go back to point 1 until all nodes are marked as visited.

At the end of the algorithm, we can compute

$$\mathcal{H}(o_1, \dots, o_{t-1}|i) = \sum_{j=1}^{t-1} \mathcal{H}(o_j|s_{o_j}) .$$

Results for bond percolation on scale-free networks

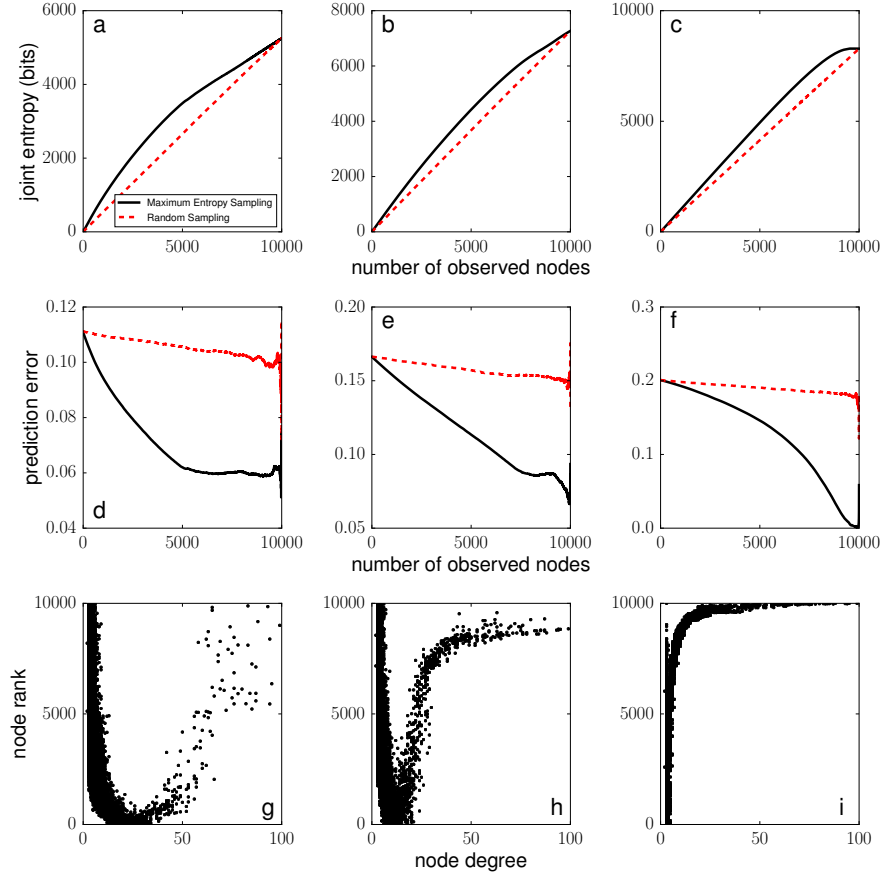


Figure A3. Maximum entropy sampling in artificial networks. We consider the bond percolation model applied to a scale-free network generated according to the uncorrelated configuration model. The degree distribution of the network is $P(k) \sim k^{-5/2}$ for $k \in [3, \sqrt{N}]$, and $P(k) = 0$, otherwise. The size of the network is $N = 10000$. The percolation threshold computed with belief propagation is $p_c = 0.0662$. Panels a, d and g refer to the case $p = 1.5p_c$, panels b, e and f to $p = 2p_c$, and panels c, f, and i to $p = 4p_c$. The description of panels a–f is similar to the one the panels of Figure 1. For random sampling, we report the results of a single realization. Inference tests are based on $T = 100$ independent simulations of the bond percolation process. In panels g, h and i every point is a node in the network. Node rank corresponds to the order in which a node is added to the set of observed nodes according to the maximum entropy sampling algorithm. This quantity is plotted against the degree of the node.

Results for the Susceptible-Infected-Susceptible (SIS) model on scale-free networks

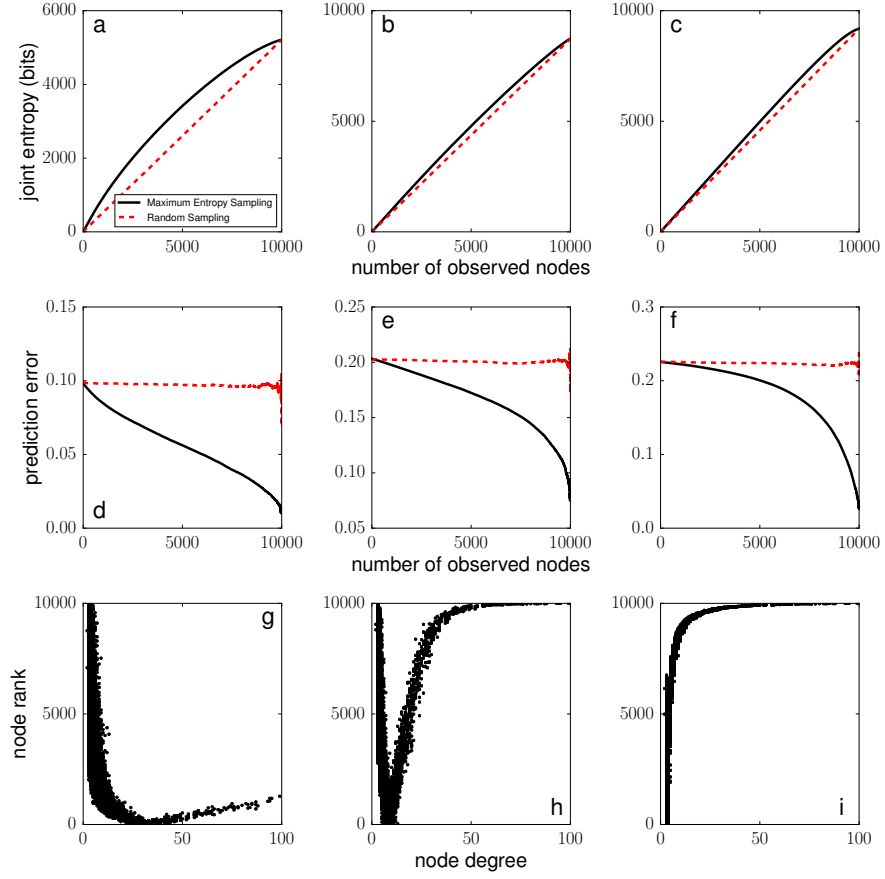


Figure A4. Maximum entropy sampling in artificial networks. We consider the SIS model applied to a scale-free network generated according to the uncorrelated configuration model. The network used here is the same as the one of considered in Fig. A3. We consider the stationary state of the SIS model for different values of the spreading rate λ (recovery rate is fixed to one). Every node i can assume only two states: $x_i = 1$ if in state I , and $x_i = 0$ if in state S . We rely on the theory developed in Ref. [27] to estimate the critical value $\lambda_c = 0.0567$ for our network. Also, we take advantage of Equation (2) of Ref. [27] to estimate the probabilities $p(x_i)$ for every node i , and $p(x_i|x_j)$ for every pair of nodes i and j . The quantities are used to compute unconditional and pairwise conditional entropies required by the maximum entropy sampling algorithm. To perform inference, we run $T = 100$ independent simulation of the SIS model. We then compare the result of a simulation with those predicted by Equation (2) of Ref. [27], where the state $\tilde{x}_{o_1}, \dots, \tilde{x}_{o_O}$ of the O observed nodes is blocked to the value observed in the simulation to infer the probability for the unobserved nodes. The error in the prediction is measured as $\epsilon^2 = \frac{1}{N-O} \sum_{i \notin O} [p(x_i = 1 | \mathbf{x}_O = \tilde{\mathbf{x}}_O) - \tilde{x}_i]^2$. The quantity is averaged over T simulations. Panels a, d and g refer to the case $\lambda = 2\lambda_c$, panels b, e and f to $\lambda = 4\lambda_c$, and panels c, f, and i to $\lambda = 8\lambda_c$. The description of panels is similar to the one the panels of Figure A3. For random sampling, we report the results of a single realization.

Independent Cascade Model on a Star-like Network

To provide further evidence on how the details of a stochastic system may affect the identification of the best set of observed nodes, we study the Independent Cascade Model (ICM) on a star-like configuration. The topology is given by a network composed of N peripheral nodes, with labels $p = 1, \dots, N$, attached to a central node with label $c = N + 1$. In the ICM model, nodes can be in three different states: S for susceptible, I for infected, and R for recovered. The rules of the dynamics are simple. At every instant of time, every infected node passes the infection to its all neighbors that are in the S state with probability λ . After that, the nodes recover, and they do longer participate in the dynamics. The process continues if new nodes have been infected, otherwise, it stops leaving us with final configuration where nodes are only in the states S or R . Here, we focus our attention on the final configuration of the system. The infection probability λ is one ingredient which we can play with. The second ingredient is the initial configuration of the system, namely $\mathbf{x}^{(0)}$, where we assume that for a generic node i we can have $x_i^{(0)} = I$ or S . In the various cases below, we compute the entropy associated with the central node or with a peripheral node. The node with larger entropy is the one initially observed according to the maximum entropy principle.

Unknown initial configuration

Let us consider the case in which we don't know anything about the initial configuration, so that every initial configuration has the same probability of appearance. The total number of possible initial configurations is 2^{N+1} , and their associated probability is $2^{-(N+1)}$.

Suppose the total number of initially infected nodes is n . The probability that one these configurations includes among the selected nodes the central node is given by

$$p(x_c^{(0)} = I | N_I = n) = \frac{\binom{N}{n-1}}{\binom{N+1}{n}} = \frac{N!}{(n-1)!(N-n+1)!} \frac{n!(N+1-n)!}{(N+1)!} = \frac{n}{N+1}.$$

The probability that the total number of nodes that are initially infected is n reads as

$$p(N_I = n) = 2^{-(N+1)} \binom{N+1}{n}. \quad (\text{A5})$$

The probability that the final state of the central node is R depends on the initial configuration. In particular, this depends on whether the central node is initially infected or not. We have

$$p(x_c = R | N_I = n, x_c^{(0)}) = \begin{cases} 1 & , \text{ if } x_c^{(0)} = I \\ 1 - (1 - \lambda)^n & , \text{ if } x_c^{(0)} = S \end{cases}$$

For a peripheral node, we have

$$p(x_p = R | N_I = n, x_c^{(0)}) = \begin{cases} \frac{n-1}{N} + \lambda \left(1 - \frac{n-1}{N}\right) & , \text{ if } x_c^{(0)} = I \\ \frac{n}{N} + \left(1 - \frac{n}{N}\right) \lambda (1 - (1 - \lambda)^n) & , \text{ if } x_c^{(0)} = S \end{cases}$$

We know that

$$p(x_c = R) = \sum_{n=0}^{N+1} p(N_I = n) \sum_{x_c^{(0)}=S,I} p(x_c^{(0)} | N_I = n) p(x_c = R | N_I = n, x_c^{(0)})$$

thus

$$p(x_c = R) = \frac{1}{2^{N+1}} \sum_{n=0}^{N+1} \binom{N+1}{n} \left\{ \frac{n}{N+1} + \left(1 - \frac{n}{N+1}\right) [1 - (1 - \lambda)^n] \right\} \quad (\text{A6})$$

For a peripheral node, we have instead

$$p(x_p = R) = \sum_{n=0}^{N+1} p(N_I = n) \sum_{x_c^{(0)}=S,I} p(x_c^{(0)} | N_I = n) p(x_p = R | N_I = n, x_c^{(0)})$$

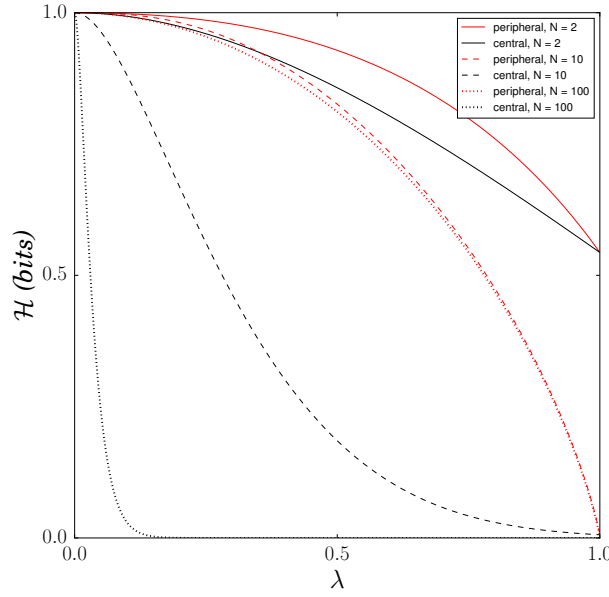


Figure A5. Entropy derived from Eqs. (A6) and (A7).

thus

$$p(x_p = R) = \frac{1}{2^{N+1}} \sum_{n=0}^{N+1} \binom{N+1}{n} \left\{ \frac{n}{N+1} \frac{n-1}{N} + \frac{n}{N+1} \lambda \left(1 - \frac{n-1}{N}\right) \right. \\ \left. \left(1 - \frac{n}{N+1}\right) \frac{n}{N} + \left(1 - \frac{n}{N+1}\right) \left(1 - \frac{n}{N}\right) \lambda [1 - (1-\lambda)^n] \right\} \quad (\text{A7})$$

Eqs. (A6) and (A7) can be finally used to compute the entropy associated with the central node or a generic peripheral node for given values of λ and N according to

$$\mathcal{H}(i) = p(x_i = R) \log_2[p(x_i = R)] + [1 - p(x_i = R)] \log_2[1 - p(x_i = R)] , \quad \text{with } i = c \text{ or } p . \quad (\text{A8})$$

Initial configuration with a fixed density of infected nodes

If we replace the probability of Eq. (A5) with the binomial distribution

$$p(N_I = n) = \binom{N+1}{n} \rho^n (1-\rho)^{N+1-n} , \quad (\text{A9})$$

every single node has a probability ρ to be initially in the I state. We can therefore replace Eq. (A6) with

$$p(x_c = R) = \sum_{n=0}^{N+1} \binom{N+1}{n} \rho^n (1-\rho)^{N+1-n} \left\{ \frac{n}{N+1} + \left(1 - \frac{n}{N+1}\right) [1 - (1-\lambda)^n] \right\} \quad (\text{A10})$$

and Eq. (A7) with

$$p(x_p = R) = \sum_{n=0}^{N+1} \binom{N+1}{n} \rho^n (1-\rho)^{N+1-n} \left\{ \frac{n}{N+1} \frac{n-1}{N} + \frac{n}{N+1} \lambda \left(1 - \frac{n-1}{N}\right) \right. \\ \left. \left(1 - \frac{n}{N+1}\right) \frac{n}{N} + \left(1 - \frac{n}{N+1}\right) \left(1 - \frac{n}{N}\right) \lambda [1 - (1-\lambda)^n] \right\} \quad (\text{A11})$$

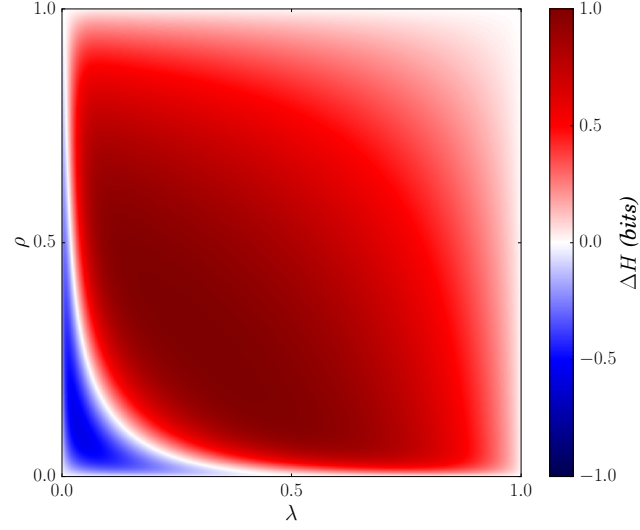


Figure A6. Difference between the entropy computed with Eq. (A11) and the one obtained with Eq. (A10) as a function of the probabilities λ and ρ . The size of the star is $N = 100$.

Initial configuration with a single infected node

Finally, let us focus only on starting configurations where the one and only one node is infected and all others are in the S state. The probability that the central node is in the final configuration in state R is

$$p(x_c = R) = \frac{1}{N+1} + \left(1 - \frac{1}{N+1}\right)\lambda \quad (\text{A12})$$

The probability that a generic peripheral node is in state R in the final configuration is

$$p(x_p = R) = \frac{1}{N+1} + \left(1 - \frac{1}{N+1}\right)\frac{1}{N}\lambda + \left(1 - \frac{1}{N+1}\right)\left(1 - \frac{1}{N}\right)\lambda^2 \quad (\text{A13})$$

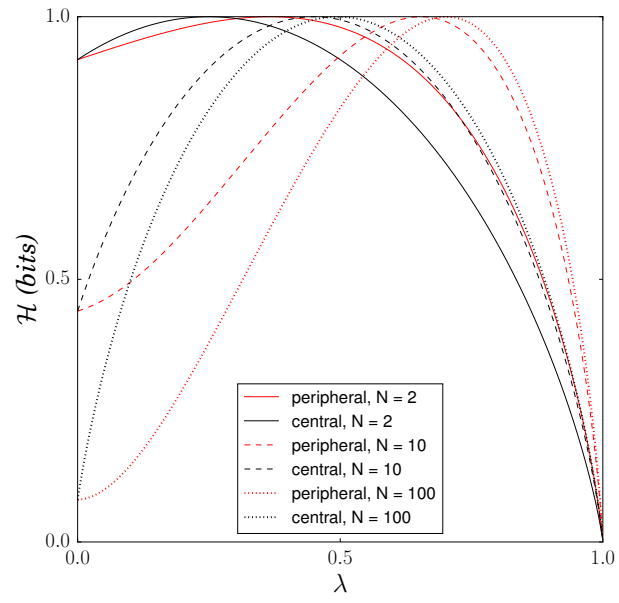


Figure A7. Entropy derived from Eqs. (A12) and (A13).